# Unsupervised Domain Adaptation with Contrastive Learning for OCT Segmentation

Alvaro Gomariz[1], Huanxiang Lu[1], Yun Yvonna Li[1], Thomas Albrecht[1],
Andreas Maunz[1], Fethallah Benmansour[1], Alessandra M. Valcarcel[2],
Jennifer Luu[2], Daniela Ferrara[2], and Orcun Goksel[3,4]

[1] F Hoffmann-La Roche AG, Basel, Switzerland
[2] Genentech Inc, California, United States
[3] Computer-assisted Applications in Medicine, ETH Zurich, Zurich, Switzerland
[4] Department of Information Technology, Uppsala University, Uppsala, Sweden

**Abstract.** Accurate segmentation of retinal fluids in 3D Optical Coherence Tomography images is key for diagnosis and personalized treatment of eye diseases. While deep learning has been successful at this task, trained supervised models often fail for images that do not resemble labeled examples, e.g. for images acquired using different devices. We hereby propose a novel semi-supervised learning framework for segmentation of volumetric images from new unlabeled domains. We jointly use supervised and contrastive learning, also introducing a contrastive pairing scheme that leverages similarity between nearby slices in 3D. In addition, we propose channel-wise aggregation as an alternative to conventional spatial-pooling aggregation for contrastive feature map projection. We evaluate our methods for domain adaptation from a (labeled) source domain to an (unlabeled) target domain, each containing images acquired with different acquisition devices. In the target domain, our method achieves a Dice coefficient 13.8% higher than SimCLR (a state-of-the-art contrastive framework), and leads to results comparable to an upper bound with supervised training in that domain. In the source domain, our model also improves the results by 5.4% Dice, by successfully leveraging information from many unlabeled images.

## 1 Introduction

Supervised learning methods, in particular UNet [20], for segmentation of retinal fluids imaged with Optical Coherence Tomography (OCT) devices have led to major advances in diagnosis, prognosis, and understanding of eye diseases [1,10,11,21,23]. However, training these supervised deep neural networks requires large amounts of labeled data, which are costly, not always feasible, and need to be repeated for each problem domain; since trained models often fail when inference data differs from labeled examples, so-called *domain-shift*, e.g. for images from a different OCT device [22]. Unsupervised domain adaptation aims to leverage information learned from a labeled data domain for applications in other domains where only unlabeled data is available. To this end, many deep

learning methods have been proposed [25], mostly using generative adversarial networks, e.g. to translate visual appearance across OCT devices [19].

Contrastive learning (CL) aims to extract informative features in a self-supervised manner by comparing (unlabeled) data pairs in a feature subspace of a network [3,5,6,7,13,14,15,18]. A widely-adopted CL framework, SimCLR [5], generates positive image pairs from the same image via image augmentations to minimize feature distances between these pairs, while maximizing their distance from augmentations of other images as negative samples. Other CL strategies aim to successfully learn without a need for negative pairs, SimSiam [7] being a representative example. CL is commonly used for pretraining models, typically using natural images such as ImageNet [9], which are then finetuned or distilled for downstream tasks, e.g. classification, detection, or segmentation [6].

Models pretrained with natural images are of limited use for medical applications, which involve images with substantially different appearances and often with 3D content, leading to a recent focus on application-specific approaches for CL pair generation in medical context [4,8]. USCL [8] minimizes the feature distance between frames of the same ultrasound video, while maximizing the distance between frames of different videos, in order to produce pretrained models for ultrasound applications. USCL also proposes a joint semi-supervised approach, which simultaneously minimizes a contrastive and supervised *classification* loss. However, to be applicable for image segmentation, this method relies on subsequent finetuning, which is potentially sub-optimal for preserving the unlabeled information for the intended task of segmentation. In fact, there exist little work on CL methods on image segmentation without finetuning.

We hereby aim to improve segmentation quality of OCT datasets with limited manual annotations, but with abundant unlabeled data. We focus on unsupervised domain adaptation, where manual annotations exist for one device (source domain), but not for another (target domain). We achieve this with the following contributions: • We introduce a semi-supervised framework for joint training of CL together with segmentation labels (Section 2.1). • We propose an augmentation strategy that leverages expected similarity between nearby slices in 3D (Section 2.2). • We introduce a new CL projection head (Section 2.3) that aggregates features without losing spatial context, which produces results superior to the conventional spatial pooling strategy. Our contributions are tested on two large clinical datasets collected in trials using different OCT imaging devices.

## 2  Methods

### 2.1  Simultaneously learning from labeled and unlabeled data

As the segmentation backbone, we utilize the proven UNet architecture [20], which can be modeled as $F(\cdot)$ processing an image $x$ to produce a segmentation map $p = F(x)$ to approximate an (expert-annotated) ground truth segmentation $y$. In the supervised setting, $F$ is learned by minimizing a supervised loss

$\mathcal{L}_{\text{sup}}$, which is for us the logarithmic Dice loss of labeled data in a domain $D$:

$$\mathcal{L}_{\text{sup}} = - \sum_{(p_i, y_i) \in D} \log \frac{2 \sum_{j \in \text{pixels}} y_i^j p_i^j}{\epsilon + \sum_{j \in \text{pixels}} (y_i^j + p_i^j)} \tag{1}$$

for all training images $i$ in $D$, where $\epsilon$ is a small number to avoid division by 0.

Contrastive frameworks aim to learn features $h = E(x)$ with an encoder $E(\cdot)$ without the need of manually annotated labels $y$. We herein base our methods on the SimCLR framework [5]. In order to adapt the learned features $h$ for our intended segmentation task, we replace the originally-proposed ResNet architecture for $E(\cdot)$ with the UNet encoder (illustrated in brown in Fig. 1a). A subsequent contrastive projection head $C(\cdot)$ maps the bottleneck-layer features to vector projections $z = C(h)$ on which the contrastive loss $\mathcal{L}_{\text{con}}$ is applied. This loss aims to minimize the distance between "positive" pairs of images $(x_i', x_i'')$ created from each image $x_i$ by a defined pair generator $P(\cdot)$ described further in Section 2.2 below, i.e. $P(x_i) = (x_i', x_i'')$. We employ a version of the normalized temperature-scaled cross entropy loss [18] adapted to our problem setting as:

$$L_{\text{con}}^{\text{CLR}} = \sum_{P(x_i),\ x_i \in D} \left( l(z_i', z_i'') + l(z_i'', z_i') \right) \tag{2}$$

$$l(z_i', z_i'') = - \log \frac{\exp\left( d(z_i', z_i'')/\tau \right)}{\sum_{x_i \in D} \mathbb{1}_{[k \neq i]} \exp\left( d(z_i', z_k'')/\tau \right)} \tag{3}$$

where $d(u, v) = (u \cdot v)/(||u||_2\, ||v||_2)$ and $\tau$ is the temperature scaling parameter.

In SimSiam, a learnable predictor $Q(\cdot)$ is applied on one projection to predict the other:

$$L_{\text{con}}^{\text{Siam}} = - \sum_{x_i \in D} \left( d\big(Q(z_i'), z_i''\big) + d\big(Q(z_i''), z_i'\big) \right) \tag{4}$$

where the gradients from the second projection pairs are prevented from backpropagating for network weight updates (*stopgrad*).

We adapt the USCL joint training strategy, which was proposed for US video classification, to our segmentation task on 3D images by combining $\mathcal{L}_{\text{sup}}$ and $\mathcal{L}_{\text{con}}$ in a semi-supervised framework illustrated in Fig. 1a. Considering a source domain $D^{\text{s}}$ and a target domain $D^{\text{t}}$, total loss $\mathcal{L}$ is calculated as follows:

$$\mathcal{L} = \frac{1}{2} \left( \underset{x \in D^{\text{s}}}{\mathcal{L}_{\text{con}}} + \underset{x \in D^{\text{t}}}{\mathcal{L}_{\text{con}}} \right) + \lambda \underset{(x,y) \in D^{\text{s}}}{\mathcal{L}_{\text{sup}}} \tag{5}$$

## 2.2 Pair generation strategy

Generation of pairs for the contrastive loss is key for successful self-supervised learning. We herein propose and compare different pair generation functions $P(\cdot)$ for volumetric OCT images, as illustrated in Fig. 1b.
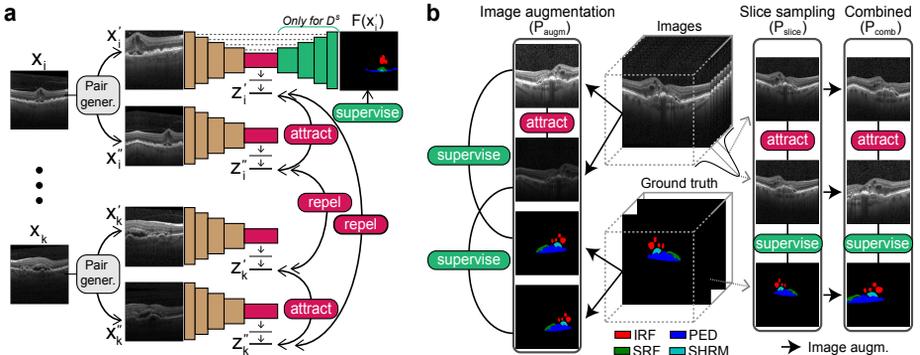
**Fig. 1.** Illustration of our CL methods. (a) Semi-supervised contrastive learning framework for unsupervised domain adaptation. Note that the *repel* modules do not apply to SimSiam. (b) Proposed pair generation methods for contrastive learning on 3D images.

We denote by $P_{\mathrm{augm}}$ an OCT adaptation of the pair formation typically employed for natural images (e.g., in SimCLR and SimSiam). Here, labeled slices in $D^{\mathrm{s}}$ and random slices in $D^{\mathrm{t}}$ are augmented with horizontal flipping ($p = 0.5$), horizontal and vertical translation (within 25% of the image size), zoom in (up to 50%), and color distortion (brightness up to 60% and jittering up to 20%). For color augmentation, images are transformed to RGB, and then back to grayscale.

We propose $P_{\mathrm{slice}}$ that leverages the coherence of nearby slices in a 3D volume for CL. Here, $x_i' = x_i$ for a slice index $b_i'$ in 3D. Then, $x_i''$ is a slice from the same volume with the (rounded) slide index $b_i'' \sim \phi(b_i', \sigma)$, where $\phi$ is a Gaussian distribution centered on $b_i'$, with standard deviation $\sigma$ as a hyperparameter. Combining the two pairing strategies yields $P_{\mathrm{comb}}$ where $P_{\mathrm{slice}}$ is used first and the augmentations in $P_{\mathrm{augm}}$ are then applied on the selected slices.

### 2.3   Projection heads to extract features for image segmentation

A projection head $C(\cdot)$ is formed by an aggregation function $\rho^{\mathrm{agg}}$ that aggregates features $h$ to form a vector, which is then processed by a multilayer perceptron $\rho^{\mathrm{MLP}}$ to create projection $z$. Typical contrastive learning frameworks, e.g. SimCLR and SimSiam, use a projection (denoted herein by $C_{\mathrm{pool}}$) where $\rho_{\mathrm{pool}}^{\mathrm{agg}} : \mathbb{R}^{w \times h \times c} \to \mathbb{R}^{1 \times 1 \times c}$ is a global pooling operation on the width $w$, height $h$, and channels $c$ of the input features. Such projection $C_{\mathrm{pool}}$ may be suboptimal for learning representations to effectively leverage segmentation information, as backpropagation from $\mathcal{L}_{\mathrm{con}}$ would lose the spatial context. Instead we propose $C_{\mathrm{ch}}$, for which $\rho_{\mathrm{ch}}^{\mathrm{agg}} : \mathbb{R}^{w \times h \times c} \to \mathbb{R}^{w \times h \times 1}$ is a $1 \times 1 \times 1$ convolutional layer that learns how to aggregate layers, so the spatial context is preserved.

# 3   Experiments and Results

**Dataset.**   We employ two large OCT datasets from clinical trials on patients with neovascular age-related macular degeneration. Images acquired using a *Spectralis* (Heidelberg Engineering) imaging device have $512 \times 496 \times 49$ or $768 \times 496 \times 19$ voxels, with a resolution of $10 \times 4 \times 111$ or $5 \times 4 \times 221$ $\mu$m/voxel, respectively. These were acquired as part of the phase-2 AVENUE trial (NCT 02484690). Images acquired as part of another study, phase-3 HARBOR trial (NCT 00891735), were acquired with a *Cirrus* HD-OCT III (Carl Zeiss Meditec) imaging device, which produces scans with $512 \times 128 \times 1024$ voxels and a resolution of $11.7 \times 47.2 \times 2.0$ $\mu$m/voxel. All slices (B-scans) from the two different devices are resampled to $512 \times 512$ pixels with roughly the same resolution of $10 \times 4$ $\mu$m/pixel. Select B-scans from Spectralis were manually annotated for fluid regions of potential diagnostic value: intraretinal fluid (IRF), subretinal fluid (SRF), pigment epithelial detachment (PED), and subretinal hyperreflective material (SHRM). More details on these datasets and the annotation protocol can be found in [17]. In our experiments, we use all training data from Spectralis as source domain $D^{\mathrm{s}}$, and unlabeled images from Cirrus as target domain $D^{\mathrm{t}}$. Labeled data from Cirrus is only used for the training of an *UpperBound* model for $D^{\mathrm{t}}$. Data stratification used in our evaluations is detailed in the supplementary Table S1.

**Implementation.**   Adam optimizer [16] was used in all models, with a learning rate of $10^{-3}$. Dropout with $p = 0.5$ is applied before and after each convolutional block in the lowest UNet resolution level, as well as after the convolutions in the two subsequent resolution levels of the decoder. Group normalization [26] with 4 groups is used after each convolutional layer. After the aggregation function $\phi$ in $C(\cdot)$, two fully-connected layers are used with 128 units each, where the first one uses ReLU activation. We heuristically set $\lambda = 20$ and the standard deviation of $\phi$ for $P_{\mathrm{slice}}$ as $\sigma = 0.25\,\mu$m, which is the range for which we observe roughly similar features across slices. Implementation is in Tensorflow 2.7, ran on an NVIDIA V100 GPU.

**Metrics.**   We segment individual slices with 2D UNet, since (1) only some slices were annotated in OCT volumes; and (2) this enables our slice-contrasting scheme. Model performance was evaluated also slice-based, using the Dice coefficient and Unnormalized Volume Dissimilarity (UVD) on 2D slices. The latter measures the extent of total segmentation error (FP+FN) in each slice and is more robust to FP on B-scans with small annotated regions for individual classes. Averaging metrics across classes with a large variation may lead to bias. Thus, we first normalize each per-slice metric ($m_i^c$) for method $i$ and class $c$ by its class Baseline ($m_{\mathrm{bas}}^c$), and then average these over all $c$ and images on the test set. All models with supervision were trained for 200 epochs, and the model at the epoch with the highest average Dice coefficient across classes on the validation set was selected for evaluation on a holdout test set.

**Table 1.** Evaluation on target domain $D^t$ and source domain $D^s$ across all classes, relative to Baseline (rel) and absolute values (abs), in red when metrics are inferior, and in bold for the best performance (excluding UpperBound). Supervised methods use labels from the domain in brackets. Dice is shown as %, and UVD as $\mu m^3 \text{x} 10^2$.

| Approach | Methods | Domain $D^t$ | | | | Domain $D^s$ | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | Dice rel | (abs) | UVD rel | (abs) | Dice rel | (abs) | UVD rel | (abs) |
| Supervised | UpperBound[$D^t$] | 29.32 | 63.88 | −8.93 | 8.67 | - | - | - | - |
| | Baseline[$D^s$] | 0.00 | 34.57 | 0.00 | 17.60 | 0.00 | 67.36 | 0.00 | 5.80 |
| Adversarial | CycleGAN [24] | −6.53 | 28.04 | 2.51 | 20.10 | −35.13 | 32.23 | 7.62 | 13.42 |
| | DAN [2] | 17.93 | 52.49 | −5.25 | 12.34 | −0.51 | 66.85 | 0.02 | 5.82 |
| Finetuning (CL → supervision) | SimCLR [5] | 14.01 | 48.58 | −4.24 | 13.36 | −3.48 | 63.88 | 0.48 | 6.28 |
| | SimSiam [7] | 11.41 | 45.97 | −2.39 | 15.21 | 0.40 | 67.75 | 0.19 | 6.00 |
| Joint (CL + supervision) | SegCLR($P_{augm},C_{pool}$) | 23.22 | 57.78 | −5.91 | 11.68 | −0.65 | 66.71 | 0.00 | 5.80 |
| | SegSiam($P_{augm},C_{pool}$) | −21.90 | 12.67 | 48.09 | 65.69 | −46.58 | 20.78 | 48.31 | 54.11 |
| | SegCLR($P_{slice},C_{pool}$) | 6.14 | 40.71 | −2.81 | 14.79 | −15.14 | 52.22 | 2.26 | 8.06 |
| | SegCLR($P_{comb},C_{pool}$) | 27.21 | 61.77 | −6.25 | 11.34 | 1.48 | 68.83 | 0.18 | 5.98 |
| | SegCLR($P_{comb},C_{ch}$) | **27.77** | **62.33** | **−6.71** | **10.88** | **1.93** | **69.28** | **−0.09** | **5.71** |

### 3.1 Evaluation on the unlabeled target domain

We first evaluate our proposed methods in the desired setting of unsupervised domain adaptation; i.e. models trained on $(x, y) \in D^s$ and $x \in D^t$ are evaluated on $y \in D^t$. Note that, although unlabeled for training, $D^t$ has some ground truth annotations in the test set to enable its evaluation (see Table S1). In Table 1 and Table S2, *UpperBound* results for a supervised model trained on labeled data from the target domain are also reported for comparison. This labeled data, used here as a reference, is ablated for all other models. A supervised UNet model, *Baseline*, was trained only on the source domain $D^s$. Its poor performance on $D^t$ confirms that the two domains indeed differ from supervised learning perspective. **Adversarial** approaches are included as state-of-the-art baselines for unsupervised domain adaptation. CycleGAN [24] is adapted to our UNet using entire slices. Training converged with meaningful translated images from $D^t$ to $D^s$, on which we run the pretrained UNet. Domain Adversarial Neural Network (DANN) includes a gradient reversal layer [12] with the design in [2] for segmentation. While DANN performs better than Baseline on $D^t$, CycleGAN is inferior. Our latter observation is contrary to that reported in [24], which is likely due to our Baseline being much superior to that of [24] (with a reported Dice of near zero). **Finetuning.** Learning representations of $D^t$ with SimCLR and SimSiam with subsequent finetuning on $D^s$ shows a clear improvement over Baseline for all classes, confirming that these CL strategies are also valid when adapted to our OCT dataset. SimCLR produces better results than SimSiam, suggesting that the use of negative pairs helps in learning better representations in our case. **Joint training** using the SimCLR framework and our above changes for a supervised loss for segmentation is herein called *SegCLR* (*SegSiam* for the SimSiam equivalent), which increases the number of parameters merely by 6.85% (7.33% for SegSiam). SegCLR($P_{augm},C_{pool}$) shows an overall improvement over finetuning. This is not the case for SegSiam($P_{augm},C_{pool}$), which suggests that the lack of negative pairs makes it difficult to simultaneously optimize $\mathcal{L}_{sup}$ and $\mathcal{L}_{con}$;
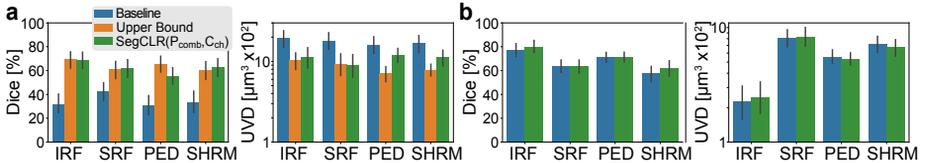
**Fig. 2.** Evaluation of models for the different classes on (a) target domain $D^{\mathrm{t}}$ and (b) source domain $D^{\mathrm{s}}$. Black bars denote 95% confidence intervals.

e.g. minimizing $\mathcal{L}_{\mathrm{con}}$ for only positive pairs may learn only simplistic features, which then would prevent $\mathcal{L}_{\mathrm{sup}}$ from improving features for segmentation.

**Pair generation.** $P_{\mathrm{slice}}$ alone produces poorer results compared to $P_{\mathrm{augm}}$ alone, indicating that merely contrasting nearby slices does not facilitate extracting features useful for segmentation. Nevertheless, by applying both pair generation methods together, i.e. with $P_{\mathrm{comb}}$, Dice and UVD results are overall superior to all the results above. This indicates that pairing nearby slices in our 3D images is a good complement to the typical image augmentation strategies.

**Projections.** We change the typical $C_{\mathrm{pool}}$ head with our proposed $C_{\mathrm{ch}}$ designed specifically for the segmentation task, which adds a mere 0.03% more parameters. While for IRF and PED (Table S2) this performs worse than SegCLR($P_{\mathrm{comb}}$,$C_{\mathrm{pool}}$), the Dice and UVD metrics averaged across classes are overall the best for SegCLR($P_{\mathrm{comb}}$,$C_{\mathrm{ch}}$), notably even surpassing the UpperBound in some cases (Fig. 2a). Hence, our proposed model could replace the UpperBound if and when no training data is available in the target domain, and in doing so only compromising the performance for PED (Fig. 2a).

### 3.2 Evaluation on the labeled source domain

Herein we test the retention of segmentation information for the original source domain $D^{\mathrm{s}}$, as shown in Table 1 (right-most column) and Fig. 2b. As expected, Baseline produces better results on $D^{\mathrm{s}}$ than on $D^{\mathrm{t}}$, since it is evaluated in the same domain in which it was supervised. For finetuning, contrary to its relative performance on $D^{\mathrm{t}}$, for $D^{\mathrm{s}}$ SimSiam produces better results than SimCLR. A reason could be SimSiam's use of only positive pairs leading to distinct features for each domain, which are later finetuned relatively more easily with segmentation supervision on $D^{\mathrm{s}}$. Further observations on $D^{\mathrm{s}}$ corroborate their above-discussed counterparts for $D^{\mathrm{t}}$; i.e. SegSiam fails; $P_{\mathrm{slice}}$ alone performs worse than $P_{\mathrm{augm}}$ alone; and combining them as $P_{\mathrm{comb}}$ performs the best. Our proposed SegCLR($P_{\mathrm{comb}}$,$C_{\mathrm{ch}}$) model produces the best results across classes also for this source domain $D^{\mathrm{s}}$, notably even surpassing the supervised Baseline. This shows that supervised information from the labeled domain is not forgotten (e.g. as a trade-off when learning from the unlabeled domain), but it is rather enhanced with the unlabeled data, despite the latter being from a different domain.
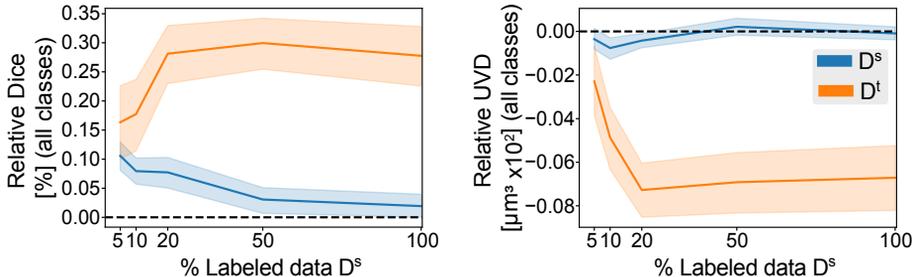
**Fig. 3.** Evaluation on $D^{\mathrm{s}}$ and $D^{\mathrm{t}}$ datasets with models trained on 5, 10, 20, 50, and 100% of labeled data from $D^{\mathrm{s}}$. Herein volume percentages are reported. Results show the proposed SegCLR($P_{\mathrm{comb}}$,$C_{\mathrm{ch}}$) relative to Baseline with same % of $D^{\mathrm{s}}$ labeled data.

### 3.3   Ablations on amount of labeled data

We study below the effect that the amount of labeled data in $D^{\mathrm{s}}$ has on the performance of our semi-supervised learning framework. To this end, we randomly ablate parts of the training data in $D^{\mathrm{s}}$. The validation set was fixed to avoid any bias on model selection. Results in Fig. 3 indicate that adding more labeled data from $D^{\mathrm{s}}$ in the training of our model has overall a positive effect on its effectiveness for segmentation of the target domain $D^{\mathrm{t}}$. This is likely because $\mathcal{L}_{\mathrm{con}}$ can adapt segmentation features to the $D^{\mathrm{t}}$ space only when these features are learned robustly with more labeled data, based on which $\mathcal{L}_{\mathrm{sup}}$ can be minimized. The trend is somewhat the opposite for $D^{\mathrm{s}}$: For the low data regime, $\mathcal{L}_{\mathrm{con}}$ seems to help with feature extraction, even though the information comes from a different domain. However, as the amount of labeled data increases and $\mathcal{L}_{\mathrm{sup}}$ is exposed to enough data from the source domain, any contrastive information contribution from a different unlabeled domain becomes relatively insignificant.

### 3.4   Segmentation results compared to inter-grader variability

Manual annotation of retinal fluids is challenging, leading to large variability in segmentation metrics even among human experts. We herein compare our proposed SegCLR($P_{\mathrm{comb}}$,$C_{\mathrm{ch}}$) to inter-grader discrepancies. We employ a set of 44 OCT volumes, each fully annotated independently by 4 different graders. These annotations are drawn from the same target domain $D^{\mathrm{t}}$ but come from a different clinical study than the dataset used in training, so a direct comparison is not possible. We evaluated segmentation metrics for graders by comparing them with one another. We deem our method within inter-grader variability when its metric for a class and image, with respect to any grader, is better than that of at least one human inter-grader metric (variation). Across images and classes, SegCLR($P_{\mathrm{comb}}$,$C_{\mathrm{ch}}$) performs within such inter-grader variability in 65.34% and 48.30% of cases based on Dice and UVD, respectively.

## 4    Conclusions

Unsupervised domain adaptation for segmentation has been typically approached as finetuning on features learned via self-supervision from classification tasks. We propose herein a segmentation approach that is jointly supervised with existing data while being self-supervised with abundant unlabeled examples from a previously unseen domain. With our proposed slice-based pairing and channel-wise aggregation for contrastive projections, our model successfully adapts supervised labeled-domain info to an unlabeled domain, surpassing previous state-of-the-art adversarial methods and even approaching the performance of an upper bound. We also improve the results in the original labeled domain by leveraging the unsupervised (contrastive) info. These contributions will help reduce manual annotation efforts for segmentation of 3D volumes in new data domains.

## References

1. Bogunovic, H., Venhuizen, F., Klimscha, S., Apostolopoulos, S., Bab-Hadiashar, A., et al.: RETOUCH: The retinal OCT fluid detection and segmentation benchmark and challenge. IEEE Transactions on Medical Imaging **38**(8), 1858–1874 (2019)
2. Bolte, J.A., et al.: Unsupervised domain adaptation to improve image segmentation quality both in the source and target domain. In: IEEE Conf on Computer Vision and Pattern Recognition (CVPR) Workshops. pp. 1404–1413 (2019)
3. Caron, M., Touvron, H., Misra, I., Jégou, H., Mairal, J., Bojanowski, P., Joulin, A.: Emerging properties in self-supervised vision transformers. In: IEEE Int Conf on Computer Vision (ICCV). pp. 9650–9660 (2021)
4. Chaitanya, K., Erdil, E., Karani, N., Konukoglu, E.: Contrastive learning of global and local features for medical image segmentation with limited annotations. In: Advances in Neural Inf Proc Systems (NeurIPS). vol. 33, pp. 12546–12558 (2020)
5. Chen, T., Kornblith, S., Norouzi, M., Hinton, G.: A simple framework for contrastive learning of visual representations. In: Int Conf on Machine Learning (ICML). pp. 1597–1607 (2020)
6. Chen, T., Kornblith, S., Swersky, K., Norouzi, M., Hinton, G.E.: Big self-supervised models are strong semi-supervised learners. In: Advances in Neural Information Processing Systems (NeurIPS). pp. 22243–22255 (2020)
7. Chen, X., He, K.: Exploring simple siamese representation learning. In: IEEE Conf on Computer Vision and Pattern Recognition (CVPR). pp. 15750–15758 (2021)
8. Chen, Y., Zhang, C., Liu, L., Feng, C., Dong, C., Luo, Y., Wan, X.: USCL: Pretraining deep ultrasound image diagnosis model through video contrastive representation learning. In: Int Conf on Medical Image Computing and Computer-Assisted Intervention (MICCAI). pp. 627–637 (2021)
9. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: ImageNet: A large-scale hierarchical image database. In: IEEE Conf on Computer Vision and Pattern Recognition (CVPR). pp. 248–255 (2009)
10. Fauw, J.D., Ledsam, J.R., Romera-Paredes, B., Nikolov, S., Tomasev, N., et al.: Clinically applicable deep learning for diagnosis and referral in retinal disease. Nature Medicine **24**(9), 1342–1350 (2018)
11. Fujimoto, J., Swanson, E.: The development, commercialization, and impact of optical coherence tomography. Investigative Ophthalmology & Visual Sci **57**(9) (2016)

12. Ganin, Y., Lempitsky, V.: Unsupervised domain adaptation by backpropagation. In: Int Conf on Machine Learning (ICML). pp. 1180–1189 (2015)
13. Grill, J.B., Strub, F., Altché, F., Tallec, C., Richemond, P., Buchatskaya, E., Doersch, C., Avila Pires, B., Guo, Z., Gheshlaghi Azar, M., et al.: Bootstrap your own latent-a new approach to self-supervised learning. In: Advances in Neural Information Processing Systems (NeurIPS). vol. 33, pp. 21271–21284 (2020)
14. He, K., Fan, H., Wu, Y., Xie, S., Girshick, R.: Momentum contrast for unsupervised visual representation learning. In: IEEE Conf on Computer Vision and Pattern Recognition (CVPR). pp. 9729–9738 (2020)
15. Khosla, P., Teterwak, P., Wang, C., Sarna, A., Tian, Y., Isola, P., Maschinot, A., Liu, C., Krishnan, D.: Supervised contrastive learning. In: Advances in Neural Information Processing Systems (NeurIPS). vol. 33, pp. 18661–18673 (2020)
16. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. In: Int Conf on Learning Representations (ICLR) (2015)
17. Maunz, A., Benmansour, F., Li, Y., Albrecht, T., Zhang, Y.P., Arcadu, F., Zheng, Y., Madhusudhan, S., Sahni, J.: Accuracy of a machine-learning algorithm for detecting and classifying choroidal neovascularization on spectral-domain optical coherence tomography. Journal of Personalized Medicine **11**(6), 524 (2021)
18. Oord, A.v.d., Li, Y., Vinyals, O.: Representation learning with contrastive predictive coding. arXiv preprint arXiv:1807.03748 (2018)
19. Ren, M., Dey, N., Fishbaugh, J., Gerig, G.: Segmentation-Renormalized deep feature modulation for unpaired image harmonization. IEEE Transactions on Medical Imaging **40**(6), 1519–1530 (2021)
20. Ronneberger, O., Fischer, P., Brox, T.: U-net: Convolutional networks for biomedical image segmentation. In: Int Conf on Medical Image Computing and Computer-Assisted Intervention (MICCAI). pp. 234–241 (2015)
21. Sahni, J.N., Maunz, A., Arcadu, F., Zhang_Schaerer, Y.P., Li, Y., Albrecht, T., Thalhammer, A., Benmansour, F.: A machine learning approach to predict response to anti-vegf treatment in patients with neovascular age-related macular degeneration using sd-oct. Investigative Ophthalmology & Visual Science **60**(11), PB094–PB094 (2019)
22. Schlegl, T., Waldstein, S.M., Bogunovic, H., Endstraßer, F., Sadeghipour, A., Philip, A.M., Podkowinski, D., Gerendas, B.S., Langs, G., Schmidt-Erfurth, U.: Fully automated detection and quantification of macular fluid in OCT using deep learning. Ophthalmology **125**(4), 549–558 (2018)
23. Schmidt-Erfurth, U., Waldstein, S.M.: A paradigm shift in imaging biomarkers in neovascular age-related macular degeneration. Progress in Retinal and Eye Research **50**, 1–24 (2016)
24. Seeböck, P., et al.: Using CycleGANs for effectively reducing image variability across oct devices and improving retinal fluid segmentation. In: IEEE Int Symp on Biomedical Imaging (ISBI). pp. 605–609 (2019)
25. Wang, M., Deng, W.: Deep visual domain adaptation: A survey. Neurocomputing **312**, 135–153 (2018)
26. Wu, Y., He, K.: Group normalization. In: European Conference on Computer Vision (ECCV). pp. 3–19 (2018)
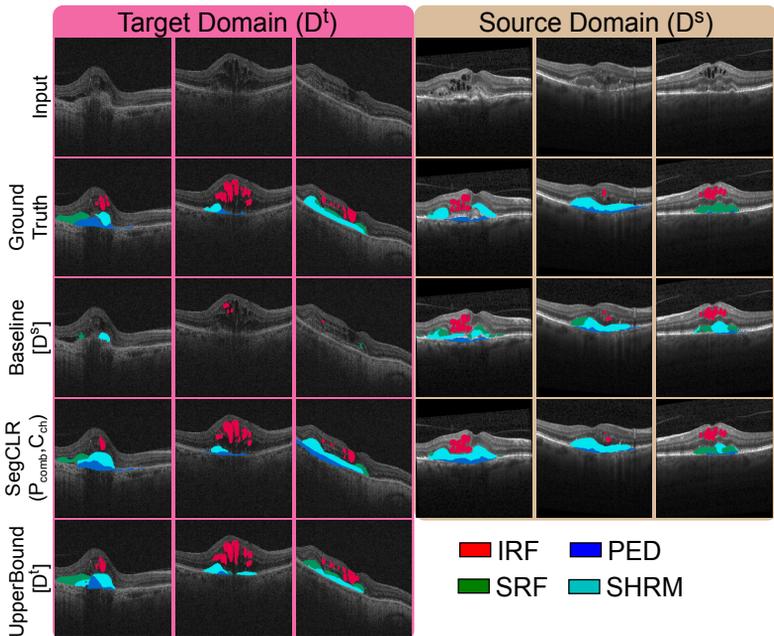
# Supplementary Material



**Fig. S1.** Qualitative assessment of segmentation on $D^{\mathrm{t}}$ and $D^{\mathrm{s}}$ examples.

**Table S1.** Datasets employed for the training and evaluation of models. Labeled data for training is displayed as #training+#validation. Volumes from the Spectralis device were used as both labeled and unlabeled data, i.e. the annotated B-scans were used for $\mathcal{L}_{\mathrm{sup}}$, while all slices were available as unlabeled data for $\mathcal{L}_{\mathrm{con}}$. Labeled training data for Cirrus (denoted in parantheses) is used only for training UpperBound.

| Domain | Device | Training | | | | Testing | |
| --- | --- | --- | --- | --- | --- | --- | --- |
| | | Labeled | | Unlabeled | | Labeled | |
| | | B-scans | Volumes | B-scans | Volumes | B-scans | Volumes |
| $D^{\mathrm{s}}$ | Spectralis | 1363+243 | 234+41 | 11 466 | 275 | 163 | 28 |
| $D^{\mathrm{t}}$ | Cirrus | (735+125) | (122+21) | 6.8 million | 53 197 | 99 | 17 |

**Table S2.** Evaluation of models on target domain $D^{\mathrm{t}}$ for all 4 annotated classes. This table corresponds to the results in Fig. 2a. Numbers in bold show the best performance for each metric and class, with a 2% tolerance, excluding UpperBound. Dice is shown as %, and UVD as $\mu m^3 \mathrm{x} 10^2$.

| Method | IRF Dice | IRF UVD | PED Dice | PED UVD | SHRM Dice | SHRM UVD | SRF Dice | SRF UVD |
|---|---|---|---|---|---|---|---|---|
| UpperBound[$D^{\mathrm{t}}$] | 69.33 | 10.44 | 65.23 | 7.09 | 60.03 | 7.82 | 60.93 | 9.32 |
| Baseline[$D^{\mathrm{s}}$] | 32.05 | 19.44 | 30.87 | 16.12 | 33.06 | 16.90 | 42.28 | 17.92 |
| SimCLR($P_{\mathrm{augm}}$,$C_{\mathrm{pool}}$) | 61.79 | 12.84 | 47.88 | 11.42 | 41.44 | 16.33 | 43.21 | 12.84 |
| SimSiam($P_{\mathrm{augm}}$,$C_{\mathrm{pool}}$) | 54.12 | 16.51 | 37.41 | 14.10 | 40.06 | 13.99 | 52.31 | 16.22 |
| SegCLR($P_{\mathrm{augm}}$,$C_{\mathrm{pool}}$) | 63.28 | 14.29 | 52.32 | 12.32 | **63.63** | **10.54** | 51.90 | 9.57 |
| SegSiam($P_{\mathrm{augm}}$,$C_{\mathrm{pool}}$) | 8.31 | 24.51 | 4.93 | 145.44 | 32.32 | 17.50 | 5.11 | 75.30 |
| SegCLR($P_{\mathrm{slice}}$,$C_{\mathrm{pool}}$) | 33.76 | 18.75 | 52.92 | 10.18 | 41.95 | 11.87 | 34.21 | 18.37 |
| SegCLR($P_{\mathrm{comb}}$,$C_{\mathrm{pool}}$) | **72.20** | **10.90** | **61.63** | **9.40** | 53.55 | 12.90 | 59.70 | 12.18 |
| SegCLR($P_{\mathrm{comb}}$,$C_{\mathrm{ch}}$) | 68.98 | 11.38 | 55.40 | 11.96 | **62.74** | 11.17 | **62.20** | **9.03** |

**Table S3.** Evaluation of models on source domain $D^{\mathrm{s}}$. This table corresponds to the results in Fig. 2b. Numbers in bold show the best performance for each metric and class, with a 2% tolerance. Dice is shown as %, and UVD as $\mu m^3 \mathrm{x} 10^2$.

| Method | IRF Dice | IRF UVD | PED Dice | PED UVD | SHRM Dice | SHRM UVD | SRF Dice | SRF UVD |
|---|---|---|---|---|---|---|---|---|
| Baseline[$D^{\mathrm{s}}$] | 77.15 | 2.28 | 71.28 | 5.55 | 57.47 | 7.17 | **63.53** | **8.20** |
| SimCLR($P_{\mathrm{augm}}$,$C_{\mathrm{pool}}$) | 73.97 | 2.55 | 67.44 | 6.57 | 53.30 | 7.36 | 60.82 | 8.67 |
| SimSiam($P_{\mathrm{augm}}$,$C_{\mathrm{pool}}$) | **79.88** | **2.24** | **72.10** | 5.70 | 56.41 | 7.31 | 62.62 | 8.73 |
| SegCLR($P_{\mathrm{augm}}$,$C_{\mathrm{pool}}$) | 75.33 | 2.47 | 71.12 | 5.18 | 58.34 | 6.90 | 62.04 | 8.64 |
| SegSiam($P_{\mathrm{augm}}$,$C_{\mathrm{pool}}$) | 31.96 | 4.33 | 4.88 | 129.08 | 44.79 | 12.69 | 1.50 | 70.35 |
| SegCLR($P_{\mathrm{slice}}$,$C_{\mathrm{pool}}$) | 52.66 | 3.10 | 63.60 | 9.58 | 46.22 | 8.94 | 46.42 | 10.63 |
| SegCLR($P_{\mathrm{comb}}$,$C_{\mathrm{pool}}$) | **79.78** | 2.41 | **73.28** | **4.92** | 57.74 | 7.99 | **64.54** | 8.59 |
| SegCLR($P_{\mathrm{comb}}$,$C_{\mathrm{ch}}$) | **80.18** | 2.48 | 71.25 | 5.33 | **62.12** | **6.74** | 63.58 | **8.29** |

**Table S4.** Balance between supervised and contrastive losses. Evaluation of the proposed SegCLR($P_{\mathrm{comb}}$,$C_{\mathrm{ch}}$) model with different values of the weighting parameter $\lambda$. Metrics here are calculated across all classes, relative to the same model with $\lambda=20$ used in all other experiments. Very low and high values (i.e., $\lambda = \{0.1, 1, 1000\}$) lead to substantially worse Dice and UVD metrics on both domains. Values closer to $\lambda=20$ (i.e., $\lambda = \{10, 100\}$) have only a minor negative effect on the segmentation metrics. Dice is shown as %, and UVD as $\mu m^3 \mathrm{x} 10^2$.

| $\lambda$ | $D^{\mathrm{t}}$ Dice | $D^{\mathrm{t}}$ UVD | $D^{\mathrm{s}}$ Dice | $D^{\mathrm{s}}$ UVD |
|---|---|---|---|---|
| 0.1 | −12.17 | 1.25 | −14.67 | 2.62 |
| 1 | −14.84 | 2.11 | −15.20 | 1.28 |
| 10 | −4.64 | 0.57 | −1.03 | 0.28 |
| 100 | −1.41 | −0.03 | −0.01 | 0.26 |
| 1000 | −8.65 | 2.01 | −2.31 | 0.31 |